

---

# Base-and-Sugar Dual-Frame Flow Matching for RNA Co-Design

---

Anonymous Authors<sup>1</sup>

## Abstract

Current frame-based RNA generators typically represent each nucleotide with a single coordinate frame. However, our representation analysis across 11,497 static chains and 31,432 multi-state relation groups shows that a single frame is insufficient to simultaneously capture base-mediated interactions and atom-level reconstruction. Motivated by this observation, we introduce **DuetRNA**, a dual-frame model for joint RNA sequence-structure generation. Each nucleotide carries two coupled SE(3) frames: a base-centered frame anchored at the glycosidic nitrogen for base mediated relationship, and a Gram-Schmidt sugar-centered frame for sugar-backbone reconstruction. DuetRNA learns their joint distribution with SE(3) flow matching and co-generates nucleotide identities in the same forward pass, avoiding post-hoc sequence recovery by an external inverse-folding model. Our experiments across two de novo generation protocols show that DuetRNA produces RNA structures with strong folding-based self-consistency. It achieves 44.00% scTM-validity under inverse-folded sequence evaluation and 38.50% Boltz-1-based / 37.83% RhoFold-based TM-validity under direct joint generation. These results demonstrate that explicitly factorizing RNA residues into base-centered and sugar-backbone-centered frames is a strong representation choice for RNA sequence-structure generation.

## 1. Introduction

RNA molecules perform diverse biological functions through their three-dimensional structures, which are shaped not only by chain connectivity but also by rich tertiary interactions such as base pairing, base stacking, and long-range

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Submitted to the 2026 Workshop on Generative and Agentic AI for Biology (ICML 2026). Do not distribute.

packing (Saenger, 1984; Leontis & Westhof, 2001). Recent progress in RNA machine learning has substantially improved *sequence-conditioned* tertiary-structure prediction, with models such as DeepFoldRNA, DRfold, RhoFold+, and NuFold learning increasingly accurate maps from sequence to structure (Pearce et al., 2022; Li et al., 2023; Shen et al., 2024; Kagaya et al., 2025). A complementary line of work has begun to address *de novo* RNA generation, including structure-first generative models such as RNA-FrameFlow and emerging joint sequence-structure generators such as RiboGen and RiboFlow (Anand et al., 2025; Rubin et al., 2025; Ma et al., 2025). These advances suggest that RNA generative modeling is becoming feasible, but they also expose a central difficulty that is specific to RNA geometry: *what is the right generative object for a nucleotide?*

For proteins, residue-level rigid-body representations are often reasonably aligned with backbone-centered organization (Chothia, 1984). However, RNA has its own characteristics: its tertiary organization is dominated by nucleobase-mediated interactions (Leontis & Westhof, 2001), whereas its sugar-phosphate backbone remains substantially more flexible (Saenger, 1984). This asymmetry creates a representation bottleneck for the generation of single-frame residues.

**An interesting observation on frame representation analysis.** To make this bottleneck concrete, we systematically analyze 7 single-frame candidates on 11K static chains and 31K multi-state relation groups. From this analysis discussed in 5.3, we find that the tested single-frame candidates do not simultaneously provide stable base-pairing and stacking relation poses and accurate sugar-backbone reconstruction. Meanwhile, we further check using dual-frame reconstruction and observe that a base and sugar frame representation better reduce atom RMSD than any single frame. Figure 1 (b) visualizes the best dual-frame candidate: a sugar Gram-Schmidt frame combined with a chemically anchored nucleobase-plane frame.

**Our contributions.** Motivated by this observation, we propose DuetRNA, a physically informed dual-frame framework for joint RNA sequence-structure generation. DuetRNA is organized around a *base-centered frame* that captures residue-level RNA organization and a *sugar-centered frame* that supports torsion-conditioned sugar-

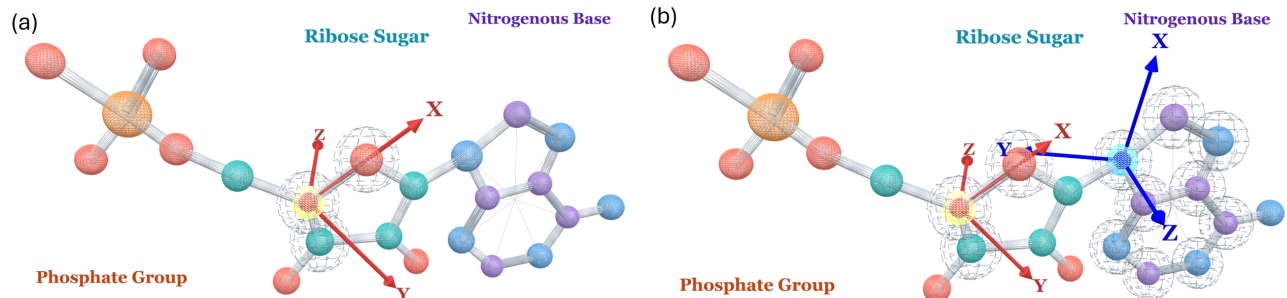


Figure 1. Comparison of RNA frame representations. Figure 1(a) builds a single a sugar Gram-Schmidt frame on  $(O4', C4', C3')$ . Figure 1(b) illustrates the dual-frame representation deployed by DuetRNA: the combination of a sugar Gram-Schmidt frame and a base frame which is chemically anchored and nucleobase-plane oriented.

backbone reconstruction. Each nucleotide is therefore represented by two coupled  $SE(3)$  frames that are jointly modeled by flow matching. We further introduce a reparameterized dual-frame constraint during training to enforce their within-residue coupling: after the network predicts the terminal base and sugar frames, we supervise the composed base-to-sugar transform derived from those predictions. This keeps the two frames chemically coupled rather than allowing them to drift as independent generated objects. With these generated frames, together with predicted sequence and intra-residue torsions, DuetRNA completes RNA atom-level reconstruction for fine-grained joint RNA sequence-structure generation.

To evaluate DuetRNA, we consider two complementary de novo generation settings. The first evaluates backbone designability through an inverse-folded sequence protocol, where generated structures are assigned candidate sequences and then independently forward-folded. The second evaluates direct joint generation, where the sequence produced by the model is directly folded by independent structure predictors, thereby testing whether the generated sequence and structure are mutually compatible without an external inverse-folding post-processing step.

Across these settings, DuetRNA produces RNA designs with strong folding-based self-consistency. It achieves 44.00% scTM-validity under inverse-folded sequence evaluation, and 38.50% Boltz-1-based / 37.83% RhoFold-based TM-validity under direct joint generation. These results indicate that a substantial fraction of the generated RNA sequence-structure pairs remain compatible under independent forward-folding evaluation. Together with our frame representation analysis, they support dual-frame factorization as an effective inductive bias for RNA generative modeling: base-centered frames capture nucleobase-mediated organization, while sugar-centered frames preserve the intra-residue geometry needed for backbone reconstruction.

## 2. Related Work

**Sequence-conditioned RNA tertiary-structure prediction.** Most recent deep-learning systems for RNA 3D modeling follow a sequence-conditioned forward-folding paradigm: given an RNA sequence  $s$ , they predict a tertiary structure  $x$ . Representative RNA-specific predictors include DeepFoldRNA, DRfold, RhoFold+, and NuFold (Pearce et al., 2022; Li et al., 2023; Shen et al., 2024; Kagaya et al., 2025), while broader biomolecular models such as RoseTTAFoldNA and AlphaFold3 extend AlphaFold-style architectures to nucleic acids and protein-nucleic-acid complexes (Baek et al., 2024; Abramson et al., 2024). These models are valuable as forward oracles in RNA design: once a candidate sequence has been proposed, they provide a practical estimate of whether it is compatible with a target fold. However, they do not directly sample RNA structures or sequence-structure pairs.

From a representation perspective, these predictors also reveal the importance of RNA-specific geometric parameterizations. DRfold learns nucleotide-wise local frames and inter-nucleotide geometric restraints, but uses a coarse-grained representation based on  $P$ ,  $C4'$ , and glycosidic  $N$  atoms, with full-atom coordinates recovered afterward (Li et al., 2023). NuFold is particularly relevant to our motivation: it adopts a flexible nucleobase-centered representation with explicit torsional modeling, and shows improved local geometry as well as accurate reconstruction of both  $C3'$ -endo and  $C2'$ -endo sugar puckers (Kagaya et al., 2025). These results suggest that base orientation and sugar-backbone conformation play distinct geometric roles in RNA modeling. At the same time, systematic benchmarks show that current prediction methods still struggle on orphan RNAs and non-Watson-Crick interactions, with local interaction fidelity remaining a major bottleneck (Bahai et al., 2024). This limitation is especially relevant for generation, where the model must construct plausible global folds together with fine-grained base-mediated interactions.

**Structure-first *de novo* RNA generation.** A complementary line of work learns a prior over RNA structural space directly, rather than only a sequence-conditioned folding map. RNA-FrameFlow introduces SE(3) flow matching for *de novo* RNA backbone generation by representing each nucleotide as a rigid-body frame and predicting the remaining backbone atoms through torsional degrees of freedom (Anand et al., 2025). Its evaluation protocol further frames generation as a design problem: sampled backbones are passed through inverse folding to obtain candidate sequences and then through forward folding to measure self-consistency (Anand et al., 2025; Joshi et al., 2025). This is a substantial step beyond forward prediction, because it learns a generative prior over RNA backbone geometry.

However, structure-first generation still separates scaffold generation from sequence compatibility. RNA-FrameFlow outputs backbone scaffolds, while compatible sequences are assigned downstream by a separate inverse-design model (Anand et al., 2025; Joshi et al., 2025). Recent generalized biopolymer design frameworks such as RFDpoly further indicate that structure-first generation for RNA, DNA, proteins, and mixed assemblies is becoming increasingly practical (Favor et al., 2025). Nevertheless, these pipelines do not directly model sequence and structure as coupled generated variables. Moreover, their generative states remain organized primarily around backbone or residue-level geometry, leaving base organization to be recovered indirectly.

**Joint sequence-structure generation and representation gap.** RiboGen moves toward direct modeling of  $p(\mathbf{s}, \mathbf{x})$  by combining continuous flow matching for geometry with discrete flow matching for sequence (Rubin et al., 2025), demonstrating that RNA co-generation is feasible. However, its frameless parameterization offers excessive modeling freedom. With no rigid-body scaffold to regularize the generation and with a training loss that does not explicitly enforce bond lengths, bond angles, or torsional preferences, the model has difficulty learning the physical constraints inside RNA, resulting in lower validity (Rubin et al., 2025). RiboFlow (Ma et al., 2025) addresses the physical-plausibility gap by retaining a single residue-level SE(3) frame (as in RNA-FrameFlow) and reconstructing atoms from backbone frame and predicted torsions, which imposes implicit geometric constraints on the output. Nevertheless, as our frame-representation analysis shows (§B), single-frame construction can not capture base-mediated interaction geometry and backbone reconstruction simultaneously.

Taken together, these create a methodological opportunity. RNA-FrameFlow establishes that SE(3) flow matching is a viable framework for RNA backbone generation (Anand et al., 2025). NuFold suggests that nucleobase-centered representations with explicit torsional structure can improve RNA local geometry (Kagaya et al., 2025). RiboGen shows

that joint RNA sequence-structure generation is possible (Rubin et al., 2025). RiboFlow (Ma et al., 2025) addresses the physical-plausibility gap left by RiboGen by retaining a single rigid frame.

Our work addresses the remaining representation gap. We formulate RNA generation around a base-sugar dual-frame state space: a base-centered frame captures inter-residue organization such as pairing, stacking, and tertiary packing, while a sugar frame captures sugar-backbone geometry and chain realization. This factorization is designed to avoid forcing high-level base organization and local backbone flexibility into a single nucleotide-level rigid body, and provides a physically informed state space for joint RNA sequence-structure generation with SE(3) flow matching.

### 3. Preliminary

#### 3.1. Rigid-Body Motions and SE(3)

**The Rotation Group SO(3).** A rotation of a rigid body is represented by a  $3 \times 3$  matrix  $R$  that preserves distances and orientations. All such matrices form the *special orthogonal group*

$$\text{SO}(3) = \{ R \in \mathbb{R}^{3 \times 3} \mid R^\top R = I, \det(R) = 1 \}. \quad (1)$$

SO(3) is a three-dimensional *Lie group* and can be viewed as a Riemannian manifold once equipped with a suitable metric. Since SO(3) is not a Euclidean vector space, operations such as interpolation, averaging, and regression must respect its manifold structure.

**The Special Euclidean Group SE(3).** A full rigid-body pose adds a translation  $t \in \mathbb{R}^3$  to the orientation. The pair belongs to the *special Euclidean group* (Barfoot, 2017)

$$\text{SE}(3) = \left\{ \begin{pmatrix} R & t \\ \mathbf{0}^\top & 1 \end{pmatrix} \mid R \in \text{SO}(3), t \in \mathbb{R}^3 \right\}. \quad (2)$$

In the text, we often denote an element compactly as  $g = (R, t)$ . Composition and inversion follow the familiar rules

$$g_1 \circ g_2 = (R_1 R_2, R_1 t_2 + t_1), g^{-1} = (R^\top, -R^\top t). \quad (3)$$

SE(3) is also a Lie group. Its tangent space at the identity, the Lie algebra  $\mathfrak{se}(3)$ , is a six-dimensional vector space. An element of  $\mathfrak{se}(3)$  is parameterized by a six-dimensional coordinate  $\xi = (\omega, v)$ , with  $\omega$  encoding the angular component and  $v$  the translational component. The exponential and logarithmic maps

$$\text{Exp} : \mathfrak{se}(3) \rightarrow \text{SE}(3), \quad \text{Log} : \text{SE}(3) \rightarrow \mathfrak{se}(3), \quad (4)$$

provide a local correspondence between the Lie algebra and the group. These maps enable local computations in the Lie algebra while representing configurations as points on SE(3).

### 3.2. Flow Matching On Riemannian Manifolds

Let  $\mathcal{M}$  be a Riemannian manifold and let  $p_{\text{data}}$  denote the unknown data distribution on  $\mathcal{M}$ . Flow matching (Lipman et al., 2023; Chen & Lipman, 2024) learns a time-dependent vector field whose induced flow transports a simple prior distribution  $p_{\text{prior}}$  on  $\mathcal{M}$  to the data distribution.

**Probability flow ODE.** A time-dependent vector field  $v_t : \mathcal{M} \rightarrow T\mathcal{M}$  defines an ordinary differential equation as

$$\frac{d}{dt} z_t = v_t(z_t), \quad z_0 \sim p_{\text{prior}}. \quad (5)$$

Here,  $T\mathcal{M}$  denotes the tangent bundle of  $\mathcal{M}$ , and  $v_t(z) \in T_z\mathcal{M}$  for each  $z \in \mathcal{M}$ . Integrating this ODE from  $t = 0$  to  $t = 1$  yields a flow  $\psi_t$  such that  $z_t = \psi_t(z_0)$ .

**Interpolation construction.** Since the marginal probability path  $p_t = (\psi_t)_\# p_{\text{prior}}$  is intractable, we adopt a simulation-free training strategy. For each data point  $z_1 \sim p_{\text{data}}$ , we draw a random noise sample  $z_0 \sim p_{\text{prior}}$  and construct an interpolation  $z_t$  between them. A natural choice on  $\mathcal{M}$  is the geodesic path:

$$z_t = \exp_{z_0}(t \log_{z_0}(z_1)), \quad t \in [0, 1], \quad (6)$$

where  $\exp_z$  and  $\log_z$  denote the Riemannian exponential and logarithmic maps at point  $z$ .

**Training via Endpoint parameterization.** We adopt an *endpoint parameterization* (Lipman et al., 2023), in which the network predicts the terminal point  $\hat{z}_1 = h_\theta(t, z_t) \in \mathcal{M}$  from interpolation point  $z_t$ . This parameterization allows the velocity to be constructed geometrically through the logarithmic map. The learning objective is then written in the tangent space at  $z_t$  as

$$\mathcal{L}(\theta) = \mathbb{E}_{t, z_0 \sim p_{\text{prior}}, z_1 \sim p_{\text{data}}} \left\| \log_{z_t}(\hat{z}_1) - \log_{z_t}(z_1) \right\|^2. \quad (7)$$

**Inference.** Once trained, the model generates samples by drawing  $z_0 \sim p_{\text{prior}}$  and numerically integrating the learned ODE from  $t = 0$  to  $t = 1$ . At each integration step, the network outputs an endpoint estimate  $\hat{z}_1$ , and the vector field is reconstructed via  $v_\theta(t, z_t) = \log_{z_t}(\hat{z}_1)/(1-t)$ . The ODE flow  $\frac{d}{dt} z_t = v_\theta(t, z_t)$ , where  $z_0 \sim p_{\text{prior}}$ , is then solved forward in time. The specific form of the tangent-space distance for our product manifold setting is given in Section 4.

## 4. Method

We propose DuetRNA, a base-and-sugar dual-frame generative framework for joint RNA sequence-structure modeling. As shown in Figure 2, each nucleotide is described by two physically meaningful rigid objects: a nucleobase-centered frame for residue-level organization and a sugar-centered frame for sugar-backbone realization. The model

learns a joint SE(3) flow over such two frames, predicts nucleotide identity and intra-residue torsional variables from coupled geometric features, and decodes the final state into an atom23 heavy-atom representation.

**Problem formulation.** Let  $\mathcal{C} = \{\text{A, U, G, C}\}$  denote the RNA nucleotide alphabet,  $s = (s_1, \dots, s_N) \in \mathcal{C}^N$  the nucleotide sequence, and  $\phi = (\phi_1, \dots, \phi_N)$  the backbone torsional variables. Our model jointly designs the structure and sequence of RNA by learning  $p_\theta(T^{\text{base}}, T^{\text{sugar}}, s, \phi)$ , which is decomposed into:

$$p_\theta(T^{\text{base}}, T^{\text{sugar}}, s, \phi) = p_\theta^{\text{frame}}(T^{\text{base}}, T^{\text{sugar}}) \cdot p_\theta^{\text{seq}}(s | T^{\text{base}}, T^{\text{sugar}}) \cdot p_\theta^{\text{loc}}(\phi | T^{\text{base}}, T^{\text{sugar}}, s). \quad (8)$$

The following subsections mirror our generative pipeline. Section 4.1 first transforms raw RNA data into a dual-frame representation. Section 4.2 estimates the joint distribution over base and sugar frames. Section 4.3 and Section 4.4 describe sequence prediction and intra-residue torsion estimation for atom23 completion. Section 4.5 defines the deterministic atom23 completion path, and Section 4.6 summarizes the training objective and model architecture.

### 4.1. Dual-Frame Representation of RNA

Systematic analysis of residue-level coordinate frames (details in Appendix B) indicates that no single frame simultaneously captures base-mediated inter-residue geometry and intra-residue backbone reconstruction geometry. Base frames preserve the inter-residue relationship more faithfully, while sugar-centered frames capture the relationship between local residue atoms required for reconstruction. To reconcile this trade-off, we adopt a dual-frame representation: a base frame  $T_i^{\text{base}}$  captures base-mediated inter-residue geometry, and a sugar frame  $T_i^{\text{sugar}}$  captures the intra-residue torsional geometry.

**The construction of sugar frame and base frame deployed in DuetRNA.** As shown in figure Figure 1, DuetRNA constructed one sugar frame and one base frame for each nucleotide. Specifically, the sugar frame  $T^{\text{sugar}}$  uses **Sugar-GS**, a three-atom Gram-Schmidt frame anchored at  $C4'$  and built from  $(O4', C4', C3')$ . The base frame  $T^{\text{base}}$  uses **Base-Plane**, a chemically anchored nucleobase-plane frame whose origin is the glycosidic connection atom ( $N9$  for purines and  $N1$  for pyrimidines) and whose orientation is defined by the fitted base plane and an in-plane chemical axis. The exact coordinate definitions of these two deployed frames are given in Appendix B.

**Dual-Frame representation for RNA.** For a chain of length

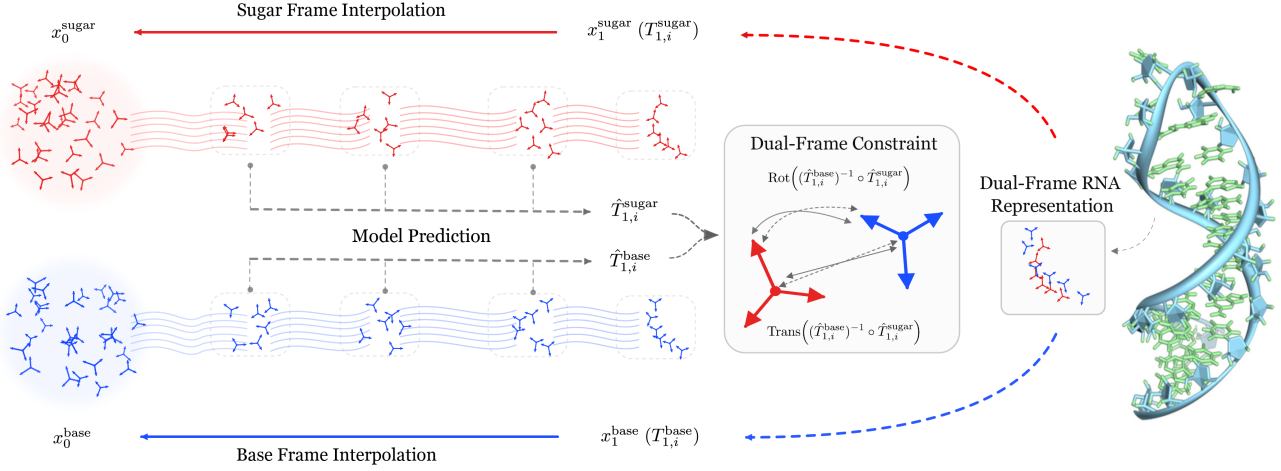


Figure 2. Overview of DuetRNA. In DuetRNA, each nucleotide is represented by two coupled SE(3) frames: a sugar frame and a base frame. The model jointly predicts terminal sugar and base frames from noisy interpolants, while a dual-frame constraint matches the predicted base-to-sugar relative translation and rotation to their ground-truth relative pose. This couples the two generated frames into a physically consistent dual-frame RNA representation for downstream sequence and atom-level structure reconstruction.

$N$ , the RNA geometric state is thus described by

$$\begin{aligned} T^{\text{base}} &= (T_1^{\text{base}}, \dots, T_N^{\text{base}}), \\ T^{\text{sugar}} &= (T_1^{\text{sugar}}, \dots, T_N^{\text{sugar}}), \end{aligned} \quad (9)$$

## 4.2. Dual-Frame Flow Matching

For an RNA molecule with  $N$  residues, the learning target for continuous geometry is the pair of frames

$$X_t = (T_t^{\text{base}}, T_t^{\text{sugar}}) \in (\text{SE}(3) \times \text{SE}(3))^N. \quad (10)$$

**Source endpoints and Dual-Frame interpolation.** For  $f \in \{\text{base}, \text{sugar}\}$ , we sample a source endpoint  $T_0^f = (R_0^f, x_0^f)$  from an unconditional source distribution with Gaussian translations and uniform rotations on  $\text{SO}(3)$  as the prior distribution  $p_0$ , and take a clean data frame  $T_1^f = (R_1^f, x_1^f)$  as the target distribution  $p_1$ . We then define the interpolation path  $T_t^f = (R_t^f, x_t^f)$  at time  $t$  as:

$$\begin{aligned} x_t^f &= (1-t)x_0^f + tx_1^f, \\ R_t^f &= \text{Exp}_{R_0^f} \left( t \text{Log}_{R_0^f} (R_1^f) \right). \end{aligned} \quad (11)$$

**Reparameterized vector field estimation.** Given  $T_t^f$  and time  $t$ , the model estimates the vector field by predicting terminal frame state  $\hat{T}_1^f$ . The implied translation and rotation vector fields from this reparameterization are

$$u_t^{\text{trans},f} = \frac{\hat{x}_1^f - x_t^f}{1-t}, \quad u_t^{\text{rot},f} = \text{Log}_{R_t^f} (\hat{R}_1^f). \quad (12)$$

**Learning objectives for vector field estimation.** The estimation objective designed for Dual-Frame flow matching

is

$$\mathcal{L}_{\text{Dual-Frame}} = \sum_{f \in \{\text{base}, \text{sugar}\}} \mathcal{L}_{\text{frame}}^f + \mathcal{L}_{\text{constraint}} \quad (13)$$

Here the rigid term is the combined SE(3) flow matching loss described in RNA-FrameFlow:

$$\begin{aligned} \mathcal{L}_{\text{frame}}^f &= \frac{1}{N} \sum_{i=1}^N \frac{\|(\hat{t}_{1,i}^f - t_{1,i}^f)\|_2^2}{(1-t)^2} \\ &+ \frac{1}{N} \sum_{i=1}^N \frac{\|\text{Log}_{R_{t,i}^f} (\hat{R}_{1,i}^f) - \text{Log}_{R_{t,i}^f} (R_{1,i}^f)\|_2^2}{(1-t)^2} \end{aligned} \quad (14)$$

To ensure the base and sugar frames preserve the within-residue relationship during the flow matching process, we further introduce the reparameterized Dual-Frame constraint term in Equation (13) which supervises terminal relative translation and rotation at the predicted terminal state ( $t = 1$ ).

$$\begin{aligned} \mathcal{L}_{\text{constraint}} &= \frac{1}{N} \sum_{i=1}^N \left\| \text{Trans} \left( (\hat{T}_{1,i}^{\text{base}})^{-1} \circ \hat{T}_{1,i}^{\text{sugar}} \right) \right. \\ &\quad \left. - \text{Trans} \left( (T_{1,i}^{\text{base}})^{-1} \circ T_{1,i}^{\text{sugar}} \right) \right\|_2^2 \\ &+ \frac{1}{N} \sum_{i=1}^N \left\| \log \left( \text{Rot} \left( (\hat{T}_{1,i}^{\text{base}})^{-1} \circ \hat{T}_{1,i}^{\text{sugar}} \right) \right) \right. \\ &\quad \left. \cdot \text{Rot} \left( (T_{1,i}^{\text{base}})^{-1} \circ T_{1,i}^{\text{sugar}} \right) \right\|_2^2. \end{aligned} \quad (15)$$

Here  $\text{Trans}(\cdot)$  and  $\text{Rot}(\cdot)$  extract the translation and rotation of an SE(3) transform, respectively.

### 4.3. Sequence Prediction

The sequence head estimates  $p_{\theta}^{\text{seq}}(s \mid T^{\text{base}}, T^{\text{sugar}})$  as a residue-wise categorical distribution. Let  $\pi_i(c)$  be the predicted probability of nucleotide  $c \in \{A, U, G, C\}$  at residue  $i$ :

$$p_{\theta}^{\text{seq}}(s \mid T^{\text{base}}, T^{\text{sugar}}) = \prod_{i=1}^N \pi_i(s_i).$$

The maximum-likelihood objective is the residue-wise negative log-likelihood

$$\mathcal{L}_{\text{Seq}} = -\frac{1}{N} \sum_{i=1}^N \log \pi_i(s_i^{\text{gt}}). \quad (16)$$

### 4.4. Intra-residue Torsion Estimation

The intra-residue torsional variables are the torsion angles required by atom23 completion:

$$p_{\theta}^{\text{loc}}(\phi \mid T^{\text{base}}, T^{\text{sugar}}, s). \quad (17)$$

Let  $K = 8$  be the number of torsion angles used by the completion map. These angles are predicted from the coupled geometric features by

$$\hat{u}_i^{\phi} = g_{\theta}^{\text{tor}}(T^{\text{base}}, T^{\text{sugar}}, s)_i \in \mathbb{R}^{K \times 2} \quad (18)$$

During training, these internal-coordinate variables are supervised both directly and through the decoded atom23 structure described in Section 4.5. The main local objective is

$$\mathcal{L}_{\text{Loc}} = \mathcal{L}_{\text{tor}} + \mathcal{L}_{\text{atom}}. \quad (19)$$

The torsion supervision uses a sine-cosine parameterization:

$$\mathcal{L}_{\text{tor}} = \frac{1}{NK} \sum_{i=1}^N \sum_{k=1}^K \left\| \hat{u}_{i,k}^{\phi} - u(\phi_{i,k}^{\text{gt}}) \right\|_2^2. \quad (20)$$

where  $u(\phi_{i,k}) = (\sin \phi_{i,k}, \cos \phi_{i,k}) \in S^1$ .

Let  $\hat{X}^{\text{atom23}}$  denote the completed atom23 structure obtained from the deterministic completion map in Section 4.5. The atom reconstruction loss is

$$\mathcal{L}_{\text{atom}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|\mathcal{A}_i|} \sum_{a \in \mathcal{A}_i} \left\| \hat{x}_{i,a} - x_{i,a}^{\text{gt}} \right\|_2^2, \quad (21)$$

where  $\mathcal{A}_i$  is the valid heavy-atom subset for residue  $i$  under the atom23 mask. Additional decoded-geometry supervision terms used in ablations are summarized in Appendix C.

### 4.5. Atom23 Completion

Atom23 denotes a fixed-width heavy-atom representation with up to 23 nucleotide-specific atom slots. Residue-specific absent atoms are masked, and hydrogens are excluded. Atom23 completion is a deterministic mapping from the predicted geometric variables to this compact heavy-atom representation:

$$\hat{X}^{\text{atom23}} = \mathcal{A}\left(\hat{T}^{\text{base}}, \hat{T}^{\text{sugar}}, \hat{\phi}, s\right), \quad (22)$$

The sugar/backbone atoms are generated from the predicted sugar frame  $\hat{T}_i^{\text{sugar}}$  together with the predicted torsions, whereas the base atoms are placed from canonical nucleotide-specific templates attached to the predicted base frame  $\hat{T}_i^{\text{base}}$ . The final output is a compact atom23 heavy-atom representation rather than a hydrogen-complete all-atom model.

### 4.6. Training Objectives And Model Architecture

**Objective.** Overall, the training loss is the sum of three components from Sections 4.2 to 4.4:

$$\mathcal{L} = \mathcal{L}_{\text{Dual-Frame}} + \lambda_{\text{Seq}} \mathcal{L}_{\text{Seq}} + \mathbf{1}[t > \tau] \lambda_{\text{Loc}} \mathcal{L}_{\text{Loc}}, \quad (23)$$

where  $\mathcal{L}_{\text{Dual-Frame}}$ ,  $\mathcal{L}_{\text{Seq}}$  and  $\mathcal{L}_{\text{Loc}}$  are described in Equations (13), (16) and (19), respectively. We use  $\tau = 0.25$  to gate  $\mathcal{L}_{\text{Loc}}$  to delay decoded-geometry supervision until the noised frames have moved beyond the earliest part of the interpolation path. The loss weights and their ablation ranges are listed in Appendix C.

**Architecture.** The network architecture contains a coupled dual-frame geometric trunk and two output heads for sequence estimation and torsion angle estimation. The trunk takes the base and sugar frames as input, refines their coupled representations across depth, and exposes final features to the estimation heads described above. A more detailed discussion is provided in Appendix A.

## 5. Experiments

### 5.1. Experimental Setup

**Datasets.** We train DuetRNA on RNA3DB (Szikszai et al., 2024), a PDB-derived RNA structure collection curated for deep-learning benchmarks with sequence- and structure-aware nonredundant splits. From the processed RNA3DB release, we retain RNA-only homomeric chains with no protein chains and modeled lengths between 40 and 150 nucleotides. RNA-FrameFlow, RiboGen, and RiboFlow were trained on RNASolo rather than RNA3DB (Adamczyk et al., 2022; Anand et al., 2025; Rubin et al., 2025; Ma et al., 2025). Since released training code and checkpoints are not available for RiboGen and RiboFlow at the time of writing, we do

Table 1. Results of de novo RNA generation (600 samples per method unless noted).

Model	Protocol	Folder	$N_T/\alpha$	Val. scTM $\uparrow$	Val. scRMSD $\uparrow$	Div. $\uparrow$	Nov. $\downarrow$
MMDiff (Anand et al., 2025)	IF	RhoFold	100/-	0.00	-	-	-
RNA-FrameFlow (paper) (Anand et al., 2025)	IF	RhoFold	50/-	41.00	-	0.610	<b>0.540</b>
RNA-FrameFlow (retrained) (Anand et al., 2025)	IF	Rhofold	50/-	26.25	23.75	<b>64.2</b>	55.4
<b>DuetRNA (Ours)</b>	IF	RhoFold	100/20	<b>44.00</b>	<b>37.00</b>	0.353	0.586
RiboGen (Rubin et al., 2025)	GS	Boltz-1	100/-	27.17	-	0.604	-
RiboGen (Rubin et al., 2025)	GS	Boltz-1	200/-	32.17	-	0.553	-
RiboGen (Rubin et al., 2025)	GS	Boltz-1	300/-	34.17	-	0.585	-
RiboFlow (Ma et al., 2025)	GS	RhoFold	50/-	34.7	-	0.577	0.562
<b>DuetRNA (Ours)</b>	GS	RhoFold	100/20	<b>37.83</b>	33.50	0.353	0.600
<b>DuetRNA (Ours)</b>	GS	Boltz-1	100/20	<b>38.50</b>	26.40	0.353	0.600

not retrain these baselines on the RNA3DB split. We therefore report their published baseline numbers as contextual references under related evaluation protocols, rather than as strict data-identical head-to-head comparisons on RNA3DB. We only retrained RNA-FrameFlow on RNA3DB for 180k iterations for comparison.

**Generation benchmark.** Following RNA-FrameFlow (Anand et al., 2025), we evaluate *de novo* generation across lengths  $L \in \{40, 50, \dots, 150\}$ , sampling 50 independent structures per length and therefore 600 structures per full run. We compare against MMDiff, RNA-FrameFlow, RiboFlow, and RiboGen under this length-controlled benchmark.

**Sequence-source protocols.** Because existing baselines differ in whether they generate nucleotide sequences, we score generated structures under two complementary protocols. **IF** denotes the inverse-folded sequence protocol: each generated backbone is assigned  $N_{\text{seq}} = 8$  candidate sequences from gRNAde (Joshi et al., 2025), folded with RhoFold (Shen et al., 2024), and scored by the best self-consistency TM-score across the eight candidates. This protocol measures backbone designability independent of the model’s own sequence channel. **GS** denotes the generated-sequence protocol: the model-generated sequence is folded directly, without an external inverse-folding model. We evaluate self-consistency TM-score with two independent structure predictors: RhoFold (Shen et al., 2024) (to align with RiboFlow’s protocol) and Boltz-1 (to align with RiboGen’s protocol) (Rubin et al., 2025; Wohlwend et al., 2024).

**Metrics.** A structure is considered valid when its self-consistency TM-score satisfies  $\text{scTM} \geq 0.45$ . We additionally report scRMSD-validity with threshold  $\text{scRMSD} \leq 4.3$  Å, qTMclust diversity, and train-set pdbTM novelty, where lower pdbTM indicates higher novelty. Unless otherwise specified, DuetRNA results use the 180K checkpoint sampled with  $N_T = 100$  and rotation exponent  $\alpha = 20$ .

## 5.2. Main Results

Table 1 reports the performance of the dual-frame state under both backbone-designability and one-shot sequence-structure generation protocols. Under IF, where sequence recovery is supplied externally by gRNAde, DuetRNA reaches 44.00% scTM-validity. This is above the published RNA-FrameFlow reference value under a related inverse-folding evaluation protocol, although the training corpora are not identical. The result indicates that the generated dual-frame geometry remains highly designable when evaluated by an established backbone-first pipeline.

The same checkpoint also remains effective when its own generated sequence is used. Under GS with RhoFold folding, DuetRNA reaches 37.83% scTM-validity without any external inverse-folding step. With Boltz-1 folding to match RiboGen’s reported protocol, DuetRNA reaches 38.50%. Together with the within-paper frame analysis, these results support the intended role of the dual-frame formulation: it improves backbone-level designability while preserving enough base-mediated information to generate compatible nucleotide identities in a single pass.

## 5.3. Frame Representation Analysis

The frame-representation analysis is separate from the de novo generation benchmark and does not use the RNA3DB training split. Its static relation pool is built from 11,497 processed RNASolo structure units (Adamczyk et al., 2022), while the multi-state pool is built from 523 single-chain NMR RNA entries comprising 31,432 relation groups. We use these pools to motivate the dual-frame factorization with two empirical findings measured on raw RNA structures before model training.

First, base-anchored coordinates are better aligned with inter-residue RNA organization. On the multi-state pool, the deployed Base-Plane frame reduces canonical-pair drift from  $13.35^\circ/1.70$  Å under Sugar-GS to  $7.78^\circ/0.61$  Å. Base-Center attains the lowest drift in this diagnostic ( $7.76^\circ/0.47$  Å), but we use Base-Plane in DuetRNA because it is chemically anchored at the glycosidic connection atom.

Table 2. Effect of the inference sampling schedule on the 120K model under IF and GS. All rows use RhoFold, 600 generated structures, and the same model weights.

$N_T$	$\alpha$	IF		GS	
		Val. scTM $\uparrow$	Div. $\uparrow$	Val. scTM $\uparrow$	Div. $\uparrow$
50	10	32.00	<b>0.585</b>	24.67	<b>0.587</b>
100	20	<b>40.67</b>	0.443	<b>31.17</b>	0.442
200	10	37.83	0.487	27.00	0.485
200	20	38.00	0.477	28.50	0.485
300	10	36.00	0.485	27.50	0.485
300	20	38.17	0.433	28.83	0.438

Second, the base and sugar frames carry complementary information rather than defining interchangeable coordinate systems. In the reconstruction-path sweep over 51,013 interaction pairs, the deployed Base-Plane + Sugar-GS dual frame reduces base-channel RMSD from 6.21 Å under Sugar-GS alone to 1.58 Å, while preserving the Sugar-GS backbone and bridge errors (1.88 Å and 2.84 Å). The Base-Center + Sugar-GS variant gives the lowest base-channel RMSD (1.09 Å), further supporting the benefit of separating base and sugar roles. The full candidate-frame catalog, static relation-separation analysis, multi-state breakdown, and reconstruction tables are reported in Appendix B.

## 5.4. Ablation Studies

### 5.4.1. INFERENCE SAMPLING SCHEDULE

Holding the 120K model, loss function, and evaluation pipeline fixed, we ablate the inference schedule by varying the integration steps  $N_T$  and rotation exponent  $\alpha$ . Table 2 shows that  $N_T = 100, \alpha = 20$  gives the strongest validity under both protocols, raising IF scTM-validity from 32.00% to 40.67% and GS scTM-validity from 24.67% to 31.17% relative to  $N_T = 50, \alpha = 10$ . Longer trajectories do not improve validity, while  $N_T = 50, \alpha = 10$  retains higher diversity.

### 5.4.2. TRAINING BUDGET

Training longer helps only up to an intermediate budget. With architecture, objective, sampler, and evaluation fixed, Table 3 shows that increasing training from 120K to 180K raises IF scTM-validity from 40.67% to 44.00% and GS scTM-validity from 31.17% to 37.83%. Further training to 190K or 200K reduces both metrics, so we use the 180K model for the main comparison.

## 6. Conclusion

In this work, we introduce DuetRNA, a dual-frame SE(3) flow matching framework for RNA sequence–structure co-generation. Motivated by a systematic analysis of seven single-frame candidates across 11,497 static chains and 31,432 multi-state relation groups, DuetRNA rep-

Table 3. Training-budget ablation under the fixed  $N_T = 100, \alpha = 20$  sampler. The 180K model is selected for the main result; further training to 190K and 200K regresses. All rows use 600 generated structures.

Training steps	IF		GS	
	Val. scTM $\uparrow$	Val. scRMSD $\uparrow$	Val. scTM $\uparrow$	Val. scRMSD $\uparrow$
120K	40.67	<b>38.83</b>	31.17	32.00
<b>180K</b>	<b>44.00</b>	37.00	<b>37.83</b>	<b>33.50</b>
190K	38.33	31.67	29.17	26.67
200K	36.67	31.33	31.00	26.00

resents each nucleotide with coupled base- and sugar-centered frames, enabling residue-level organization and intra-residue backbone reconstruction geometry to be generated within a unified flow.

Our results show that DuetRNA with this factorization yields strong folding-based self-consistency, achieving 44.00% scTM-validity under inverse-folded sequence evaluation and 38.50% Boltz-1-based / 37.83% RhoFold-based TM-validity under direct joint generation. More broadly, our analysis supports a geometric view of RNA generation in which base-centered frames capture nucleobase-mediated pairing, stacking, and long-range organization, while sugar-centered frames preserve the degrees of freedom needed for atomistic backbone reconstruction.

Current experiments focus on medium-length RNA3DB structures of 40–150 nt, and evaluating dual-frame generation across shorter fragments, longer structured RNAs, and broader benchmark settings remains an important next step. Although modeling  $(SE(3) \times SE(3))^N$  provides greater representational capacity than single-frame formulations, it also introduces additional computational cost. Future work will extend DuetRNA to conditional and functional design tasks, such as ligand-binding aptamer and riboswitch design, and develop more efficient inter-frame coupling schemes and architectural optimizations.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning, specifically within the domain of structural bioinformatics and geometric modeling. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

## References

- Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, L., Ballard, A. J., Bambrick, J., et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 630: 493–500, 2024. doi: 10.1038/s41586-024-07487-w.
- Adamczyk, B., Antczak, M., and Szachniuk, M. Rnasolo: a repository of cleaned pdb-derived rna 3d structures.

- 440 *Bioinformatics*, 38(14):3668–3670, 2022. doi: 10.1093/  
441 bioinformatics/btac386.
- 442
- 443 Anand, R., Joshi, C. K., Morehead, A., Jamasb, A. R.,  
444 Harris, C., Mathis, S. V., Didi, K., Ying, R., Hooi, B., and  
445 Liò, P. Rna-frameflow: Flow matching for de novo 3d  
446 rna backbone design. *Transactions on Machine Learning  
447 Research*, 2025.
- 448
- 449 Baek, M., McHugh, R., Anishchenko, I., Jiang, H., Baker,  
450 D., and DiMaio, F. Accurate prediction of protein–nucleic  
451 acid complexes using rosettafoldna. *Nature Methods*, 21:  
452 117–121, 2024. doi: 10.1038/s41592-023-02086-5.
- 453
- 454 Bahai, A., Kwok, C. K., Mu, Y., and Li, Y. System-  
455 atic benchmarking of deep-learning methods for tertiary  
456 rna structure prediction. *PLOS Computational Biol-*  
457 *ogy*, 20(12):e1012715, 2024. doi: 10.1371/journal.pcbi.  
458 1012715.
- 459
- 460 Barfoot, T. D. *State Estimation for Robotics*. Cambridge  
461 University Press, Cambridge, UK, 2017.
- 462
- 463 Chen, R. T. Q. and Lipman, Y. Flow matching on general  
464 geometries. In *International Conference on Learning  
465 Representations (ICLR)*, 2024.
- 466
- 467 Chothia, C. Principles that determine the structure of pro-  
468 teins. *Annual Review of Biochemistry*, 53:537–572, 1984.  
469 doi: 10.1146/annurev.bi.53.070184.002541.
- 470
- 471 Favor, A., Quijano, R., Chernova, E., Kubaney, A., Weidle,  
472 C., Esler, M. A., McHugh, L., Carr, A., Hsia, Y., Juergens,  
473 D., et al. De novo design of rna and nucleoprotein com-  
474 plexes. *bioRxiv*, 2025. doi: 10.1101/2025.10.01.679929.
- 475
- 476 Joshi, C., Jamasb, A., Viñas, R., Harris, C., Mathis, S.,  
477 Morehead, A., Anand, R., and Liò, P. gnode: Geometric  
478 deep learning for 3d rna inverse design. In *International  
479 Conference on Learning Representations*, 2025.
- 480
- 481 Kagaya, Y., Zhang, Z., Ibtihaz, N., Wang, X., Nakamura,  
482 T., Punuru, P. D., and Kihara, D. Nufold: end-to-end  
483 approach for rna tertiary structure prediction with flexible  
484 nucleobase center representation. *Nature Communica-*  
485 *tions*, 16:881, 2025. doi: 10.1038/s41467-025-56261-7.
- 486
- 487 Leontis, N. B. and Westhof, E. Geometric nomenclature  
488 and classification of RNA base pairs. *RNA*, 7(4):499–512,  
489 2001. doi: 10.1017/S1355838201002515.
- 490
- 491 Li, Y., Zhang, C., Feng, C., Pearce, R., Freddolino, P. L.,  
492 and Zhang, Y. Integrating end-to-end learning with deep  
493 geometrical potentials for ab initio rna structure pre-  
494 diction. *Nature Communications*, 14:5745, 2023. doi:  
10.1038/s41467-023-41303-9.
- Lipman, Y., Chen, R. T. Q., Ben-Hamu, H., Nickel, M.,  
and Le, M. Flow matching for generative modeling. In  
*International Conference on Learning Representations  
(ICLR)*, 2023.
- Ma, R., Zhang, Z., Wang, Z., Hua, C., Rao, J., Zhou, Z.,  
and Zheng, S. RiboFlow: Conditional De Novo RNA co-  
design via synergistic flow matching. In *The Thirty-Ninth  
Annual Conference on Neural Information Processing  
Systems*, volume 38, 2025.
- Pearce, R., Omenn, G. S., and Zhang, Y. De novo rna  
tertiary structure prediction at atomic resolution using  
geometric potentials from deep learning. *bioRxiv*, 2022.  
doi: 10.1101/2022.05.15.491755.
- Rubin, D., dos Santos Costa, A., Ponnampati, M., and Ja-  
cobson, J. Ribogen: Rna sequence and structure co-  
generation with equivariant multifold, 2025.
- Saenger, W. *Principles of Nucleic Acid Structure*. Springer-  
Verlag, New York, 1984. ISBN 0-387-90761-0. doi:  
10.1007/978-1-4612-5190-3.
- Shen, T., Hu, Z., Sun, S., Liu, D., Wong, F., Wang, J., Chen,  
J., Wang, Y., Hong, L., Xiao, J., Zheng, L., Krishnamoor-  
thi, T., King, I., Wang, S., Yin, P., Collins, J. J., and Li,  
Y. Accurate rna 3d structure prediction using a language  
model-based deep learning approach. *Nature Methods*,  
21:2287–2298, 2024. doi: 10.1038/s41592-024-02487-0.
- Szicszai, M., Magnus, M., Sanghi, S., Kadyan, S., Bouatta,  
N., and Rivas, E. Rna3db: A structurally-dissimilar  
dataset split for training and benchmarking deep learning  
models for rna structure prediction. *Journal of Molecular  
Biology*, 436(17):168552, 2024. doi: 10.1016/j.jmb.2024.  
168552.
- Wohlwend, J., Corso, G., Passaro, S., Getz, N., Reveiz,  
M., Leidal, K., Swiderski, W., Atkinson, L., Portnoi,  
T., Chinn, I., Silterra, J., Jaakkola, T., and Barzilay, R.  
Boltz-1: Democratizing biomolecular interaction mod-  
eling. *bioRxiv*, 2024. doi: 10.1101/2024.11.19.624167.  
Preprint.

## 495 A. Method Details

### 496 A.1. Model Architecture

#### 497 A.1.1. DUAL-FRAME IPA-TRANSFORMER TRUNK

498 **Input conditioning.** Each residue first receives a shared node embedding from timestep and positional encodings. Frame-  
 499 specific pose conditioning is then added for the current base and sugar poses together with a frame-type embedding that  
 500 distinguishes base tokens from sugar tokens.

501 **Shared pair representation.** The trunk constructs three pair channels: a base-view pair representation from base-node  
 502 features and base translations, a sugar-view pair representation from sugar-node features and sugar translations, and a  
 503 bridge-view pair representation from a per-residue summary of the base and sugar states together with their translation  
 504 difference. These channels are fused into a shared pair tensor used by both geometric streams.

505 **Dual geometric streams.** At block  $k$ , the trunk maintains  $T_i^{\text{base},(k)}$  and  $T_i^{\text{sugar},(k)}$  together with their corresponding feature  
 506 streams. Two invariant point attention modules process the current base and sugar poses in parallel, followed by per-stream  
 507 transitions and rigid updates

$$508 T_i^{\text{base},(k+1)} = T_i^{\text{base},(k)} \circ \Delta T_i^{\text{base},(k)}, \quad T_i^{\text{sugar},(k+1)} = T_i^{\text{sugar},(k)} \circ \Delta T_i^{\text{sugar},(k)}. \quad (24)$$

509 **Object-token mixing and dynamic edge refinement.** To enable explicit cross-object communication, the residue-wise base  
 510 and sugar tokens are periodically concatenated into a sequence of  $2L$  object tokens and processed by a joint transformer  
 511 encoder. Between blocks, residue-level summaries from the current base and sugar streams are written back into the shared  
 512 pair tensor through an edge transition, so the pair representation evolves jointly with trunk depth.

513 **Self-conditioning.** The pair stack also receives histograms computed from self-conditioned translations. During training,  
 514 with probability 0.5, a stop-gradient forward pass provides predicted base and sugar translations that are fed back into the  
 515 pair embedders. During sampling, the previous-step predictions play the same role.

### 516 Sequence and Torsion Heads.

517 Let  $H_i^{\text{base}}$  and  $H_i^{\text{sugar}}$  denote the final residue features of the two streams. The sequence head reads the final base token  
 518 together with pooled pair context from the base-view, sugar-view, and fused pair representations. The torsion head reads the  
 519 final sugar feature together with the initial sugar embedding.

### 520 A.2. Auxiliary Loss Definitions

521 **Soft Base-Template Loss.** For each candidate nucleotide  $c \in \mathcal{C}$ , let  $X_i^{\text{base}}(c)$  be the base-template atomic coordinates  
 522 of nucleotide  $c$  placed in the predicted base frame at residue  $i$ . Given predicted categorical probabilities  $\pi_i(c)$ , the  
 523 probability-weighted base-template coordinates are

$$524 \tilde{X}_i^{\text{base}} = \sum_{c \in \mathcal{C}} \pi_i(c) X_i^{\text{base}}(c). \quad (25)$$

525 The corresponding auxiliary loss is

$$526 \mathcal{L}_{\text{soft-base}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|\mathcal{B}_i^{\text{atom}}|} \sum_{a \in \mathcal{B}_i^{\text{atom}}} \left\| \tilde{x}_{i,a}^{\text{base}} - x_{i,a}^{\text{base,gt}} \right\|_2^2, \quad (26)$$

527 where  $\mathcal{B}_i^{\text{atom}}$  is the atom23 base-atom subset for residue  $i$ . This term gives the sequence probabilities structural feedback  
 528 through base-template coordinates; it does not replace the categorical cross-entropy objective.

529 **Chain and Clash Auxiliary Losses.** Let  $\mathcal{E}$  denote the valid adjacent-residue set for which both  $O3'_i$  and  $P_{i+1}$  are present.  
 530 The chain-continuity auxiliary supervises the decoded adjacent  $O3'_i$ - $P_{i+1}$  bridge:

$$531 \mathcal{L}_{\text{chain}} = \frac{1}{|\mathcal{E}|} \sum_{i \in \mathcal{E}} \left( \left\| \hat{x}_{i,O3'} - \hat{x}_{i+1,P} \right\|_2 - \left\| x_{i,O3'}^{\text{gt}} - x_{i+1,P}^{\text{gt}} \right\|_2 \right)^2. \quad (27)$$

### A.3. Inference

Inference keeps two rigid objects for every residue throughout the whole trajectory. We first sample independent source chains  $T_0^{\text{base}} \sim p_0^{\text{base}}$  and  $T_0^{\text{sugar}} \sim p_0^{\text{sugar}}$ . At step  $t$ , one coupled model evaluation takes  $(T_t^{\text{base}}, T_t^{\text{sugar}})$  as input and predicts terminal frames  $(\hat{T}_1^{\text{base}}, \hat{T}_1^{\text{sugar}})$ . Both channels are then updated in parallel:

$$x_{t+\Delta t}^{\text{base}} = x_t^{\text{base}} + \Delta t \frac{\hat{x}_1^{\text{base}} - x_t^{\text{base}}}{1-t}, \quad R_{t+\Delta t}^{\text{base}} = \text{Exp}_{R_t^{\text{base}}} \left( \alpha(t) \Delta t \text{Log}_{R_t^{\text{base}}} (\hat{R}_1^{\text{base}}) \right), \quad (28)$$

$$x_{t+\Delta t}^{\text{sugar}} = x_t^{\text{sugar}} + \Delta t \frac{\hat{x}_1^{\text{sugar}} - x_t^{\text{sugar}}}{1-t}, \quad R_{t+\Delta t}^{\text{sugar}} = \text{Exp}_{R_t^{\text{sugar}}} \left( \alpha(t) \Delta t \text{Log}_{R_t^{\text{sugar}}} (\hat{R}_1^{\text{sugar}}) \right). \quad (29)$$

If self-conditioning is enabled, the previous-step predicted base and sugar translations are fed back into the pair representation. After the final integration step, the model produces sequence probabilities and intra-residue torsion estimates from the terminal base and sugar features. Then we choose the sequence by

$$\hat{s}_i = \arg \max_{c \in \mathcal{C}} p_\theta(s_i = c | T^{\text{base}}, T^{\text{sugar}}) \quad (30)$$

and complete atom23 coordinates from the terminal pair of frames:

$$\hat{X}^{\text{atom23}} = \mathcal{A} \left( \hat{T}^{\text{base}}, \hat{T}^{\text{sugar}}, \hat{\phi}, \hat{s} \right) \quad (31)$$

to obtain the final atom23 structure.

## B. Frame Representation Analysis: Extended Results

This appendix provides the extended diagnostics underlying the frame representation analysis of §5.3. The main text reports the two core findings: base anchoring better captures inter-residue relation geometry, while a separate sugar/backbone frame is needed for local reconstruction. Here we report the candidate-frame catalog and deployed coordinate definitions, the canonicalization and continuity screen (§B.2), the static relation-separation analysis (§B.3), the per-relation multi-state breakdown including noncanonical pairs and stacking (§B.4), and the single-frame and dual-frame reconstruction sweep (§B.5).

### B.1. Candidate Frame Catalog

We audit seven residue-level single-frame candidates spanning sugar/backbone and base-anchored constructions. The display names below are used consistently in the paper; the corresponding implementation identifiers are shown in parentheses.

#### Base family.

- **Base-Plane** (`base_plane_anchor`): a chemically anchored base-plane frame. The origin is the glycosidic connection atom ( $N9$  for purines and  $N1$  for pyrimidines). The orientation is defined by the fitted nucleobase plane normal and an in-plane chemical axis ( $C4 \rightarrow C8$  for purines,  $C4 \rightarrow C2$  for pyrimidines).
- **Base-Inertial** (`base_inertial`): an inertial frame over all nucleobase heavy atoms, with residue-type-specific sign and quadrant disambiguation. The origin remains the glycosidic connection atom.
- **Base-Center** (`base_center_standard`): the base-origin control. It uses the same inertial orientation family as Base-Inertial, but moves the origin to the mass centroid of the nucleobase heavy atoms.

#### Sugar/backbone family.

- **Sugar-GS** (`baseline_geom`): the legacy three-atom Gram-Schmidt frame built from  $(O4', C4', C3')$  with origin at  $C4'$ , following the axis convention used in RNA-FrameFlow (Anand et al., 2025).
- **Sugar-Inertial-4** (`sugar4_inertial`): a four-atom mass-weighted inertial frame over  $(O4', C4', C3', C5')$  with origin at  $C4'$  and anchor-based sign disambiguation.

- **Sugar-Inertial-6** (`sugar6_inertial`): a six-atom inertial variant over  $(C1', C2', C3', C4', O4', C5')$  with origin at  $C4'$ .
- **Sugar-Ring** (`sugar5_ring_axis`): a five-membered sugar-ring frame over  $(O4', C1', C2', C3', C4')$ , with ring-plane orientation and a sugar-ring centroid origin.

### B.1.1. DEPLOYED FRAME DEFINITIONS

The deployed dual-frame model uses Sugar-GS for  $T_i^{\text{sugar}} = (R_i^{\text{sugar}}, t_i^{\text{sugar}})$  and Base-Plane for  $T_i^{\text{base}} = (R_i^{\text{base}}, t_i^{\text{base}})$ . We give their exact coordinate definitions here.

**Sugar-GS.** Let  $x_{O4'}, x_{C4'}, x_{C3'} \in \mathbb{R}^3$  be the corresponding atom coordinates. We set the origin at  $C4'$ ,

$$t_i^{\text{sugar}} = x_{C4'}, \quad e_{1,i}^{\text{sugar}} = \frac{x_{O4'} - x_{C4'}}{\|x_{O4'} - x_{C4'}\|_2}, \quad (32)$$

orthogonalize the  $C3'$  direction against  $e_{1,i}^{\text{sugar}}$ ,

$$\tilde{e}_{2,i}^{\text{sugar}} = (x_{C3'} - x_{C4'}) - \langle x_{C3'} - x_{C4'}, e_{1,i}^{\text{sugar}} \rangle e_{1,i}^{\text{sugar}}, \quad e_{2,i}^{\text{sugar}} = \frac{\tilde{e}_{2,i}^{\text{sugar}}}{\|\tilde{e}_{2,i}^{\text{sugar}}\|_2}, \quad (33)$$

and complete a right-handed frame,

$$e_{3,i}^{\text{sugar}} = e_{1,i}^{\text{sugar}} \times e_{2,i}^{\text{sugar}}, \quad R_i^{\text{sugar}} = [e_{1,i}^{\text{sugar}}, e_{2,i}^{\text{sugar}}, e_{3,i}^{\text{sugar}}]. \quad (34)$$

**Base-Plane.** Let  $\mathcal{B}_i$  denote the heavy-atom set of the nucleobase, excluding sugar atoms, and let  $a_i^{\text{link}} \in \{N9, N1\}$  be the glycosidic connection atom. The origin is

$$t_i^{\text{base}} = x_{a_i^{\text{link}}}. \quad (35)$$

We fit a plane to  $\{x_a : a \in \mathcal{B}_i\}$  and denote its unit normal by  $\tilde{n}_i$ . To resolve the normal sign, we use the residue-type-specific reference normal

$$n_i^{\text{ref}} = \begin{cases} (x_{C4} - x_{C8}) \times (x_{C2} - x_{N9}), & \text{purine,} \\ (x_{C4} - x_{C2}) \times (x_{C6} - x_{N1}), & \text{pyrimidine,} \end{cases} \quad (36)$$

and set

$$e_{3,i}^{\text{base}} = \text{sign}(\langle \tilde{n}_i, n_i^{\text{ref}} \rangle) \tilde{n}_i. \quad (37)$$

The in-plane chemical axis is

$$u_i = \begin{cases} x_{C4} - x_{C8}, & \text{purine,} \\ x_{C4} - x_{C2}, & \text{pyrimidine,} \end{cases} \quad \tilde{e}_{1,i}^{\text{base}} = u_i - \langle u_i, e_{3,i}^{\text{base}} \rangle e_{3,i}^{\text{base}}, \quad (38)$$

which gives

$$e_{1,i}^{\text{base}} = \frac{\tilde{e}_{1,i}^{\text{base}}}{\|\tilde{e}_{1,i}^{\text{base}}\|_2}, \quad e_{2,i}^{\text{base}} = e_{3,i}^{\text{base}} \times e_{1,i}^{\text{base}}, \quad R_i^{\text{base}} = [e_{1,i}^{\text{base}}, e_{2,i}^{\text{base}}, e_{3,i}^{\text{base}}]. \quad (39)$$

## B.2. Frame Canonicalization and Continuity Screening

Before comparing downstream relation or reconstruction metrics, each frame must provide a usable SE(3) supervision target. We therefore first apply a canonicalization and continuity gate. A candidate passes only if its valid rate is at least 0.95, the maximum deterministic rebuild error is at most  $10^{-5}$  degrees, and the fraction of perturbation cases with a rotation jump above  $150^\circ$  is zero. The perturbation suite includes micro sugar-pucker,  $\chi$ -rotation, phosphate-swing, and bridge-motion probes.

Sugar-Inertial-4 is therefore not used as a formal relation-stability frame in the main evidence chain. Its average continuity is not poor, but the small nonzero catastrophic-jump tail makes it unsafe as a smooth flow-matching label. We retain it only as an exploratory reconstruction diagnostic when explicitly marked as such.

Table 4. Frame canonicalization and continuity screening. The gate is strict on catastrophic rotation jumps: any nonzero  $> 150^\circ$  jump rate fails the candidate.

Frame	Valid rate	Median cont. ( $^\circ$ )	P95 cont. ( $^\circ$ )	Jump <sub>150</sub>	Gate
Sugar-GS (baseline_geom)	1.00000	1.88	5.00	0	pass
Sugar-Inertial-4 (sugar4_inertial)	0.99982	1.16	22.36	$9.8 \times 10^{-4}$	fail
Sugar-Inertial-6 (sugar6_inertial)	0.99969	0.59	5.00	0	pass
Sugar-Ring (sugar5_ring_axis)	0.99987	0.89	5.00	0	pass
Base-Plane (base_plane_anchor)	0.99584	0.00	$1.7 \times 10^{-6}$	0	pass
Base-Inertial (base_inertial)	0.99584	0.00	$1.7 \times 10^{-6}$	0	pass
Base-Center (base_center_standard)	0.99584	0.00	$1.7 \times 10^{-6}$	0	pass

### B.3. Static Relation Separation

On the full static pool of 11,497 processed RNASolo structure units, 1,765 nonredundant representatives, and 1,045 motif-rich chains, we examine which frame family more readily separates relation types and which more tightly clusters identical relation instances. We run 20 independent sub-experiments, structured as one full set and three filtered subsets (representative, motif-rich, stem-like), across five relation categories (canonical cWW pairs, noncanonical pairs, stacking, base-phosphate, and base-backbone). Canonical cWW pairs include Watson-Crick AU/UA and GC/CG pairs plus GU/UG wobble pairs. In 11 of the 20 subset-relation combinations, the frame with the best inter-group/intra-group distance ratio belongs to the sugar/backbone family while the frame with the best 1-NN retrieval accuracy belongs to the base family. This result implies that frames rooted in sugar/backbone atoms spread different relation categories farther apart, whereas frames rooted in base atoms group geometrically identical instances closer together. No single frame family achieves the best result in every channel. Table Table 5 reports the family-level split pattern.

For a residue pair  $(i, j)$ , the inter-residue relation is represented by

$$R_{ij} = R_i^\top R_j, \quad t_{ij} = R_i^\top (t_j - t_i), \quad (40)$$

where  $R_i$  and  $t_i$  are the rotation and translation of residue  $i$ .

Table 5. Family-level split pattern on the static pool. Each row reports the number of subset-relation combinations (out of 20) in which the indicated pattern holds: “sep” is the family that better separates relation types, and “ret” is the family that more tightly retrieves identical relation instances.

Channel	Split pattern (sep   ret)	# Combinations	Fraction
joint pose	sugar/backbone   base	11	0.55
joint pose	sugar/backbone   sugar/backbone	7	0.35
joint pose	base   base	1	0.05
joint pose	base   sugar/backbone	1	0.05
translation	sugar/backbone   base	13	0.65
translation	sugar/backbone   sugar/backbone	5	0.25
translation	base   base	1	0.05
translation	base   sugar/backbone	1	0.05

These counts are descriptive evidence of two distinct representational roles in RNA residue geometry, not formal statistical claims.

### B.4. Multi-State Stability: Per-Relation Breakdown

Due to the flexibility of RNA, the same molecule may adopt different conformations. A frame that separates relation types in a static snapshot does not necessarily keep the same relation stable across conformational states. NMR and other multi-model entries therefore provide a natural test bed for measuring drift across plausible states of the same molecule.

Here we extend the multi-state analysis of §5.3 to the five paper-facing frames used in the completed external run: Sugar-GS, Sugar-Inertial-6, Sugar-Ring, Base-Plane, and Base-Center. The pool consists of 523 single-chain multi-model entries (496 tier-A NMR plus 27 other multi-model structures), comprising 31,432 relation groups, 984,382 normalized relation rows, and 336 high-coverage relation keys. Variation reflects a combination of conformational flexibility and experimental uncertainty.

For each key, *orientation drift*  $\Delta R$  (median geodesic angle from the medoid rotation) and *translation drift*  $\Delta T$  (median Euclidean distance from the medoid translation) are computed under each candidate frame. A key is *base-like* in a channel if the best base-family frame outperforms all sugar/backbone-family frames.

Table 6. Aggregate multi-state stability over all 31,432 relation groups. Lower is better.

Frame	Family	$R^{\text{rel}}$ drift ( $^{\circ}$ ) $\downarrow$	$t^{\text{rel}}$ drift ( $\text{\AA}$ ) $\downarrow$
Sugar-GS (baseline_geom)	sugar/Gram-Schmidt	15.21	1.45
Sugar-Inertial-6 (sugar6_inertial)	sugar/inertial	14.60	1.40
Sugar-Ring (sugar5_ring_axis)	sugar/ring	14.44	1.31
Base-Plane (base_plane_anchor)	base/plane	12.26	1.05
Base-Center (base_center_standard)	base/center	<b>12.26</b>	<b>1.03</b>

Table 7. Per-relation multi-state stability across the five external-run frames. Lower is better. “cP” = canonical cWW pairs ( $n = 2,334$ ), “ncP” = noncanonical pairs ( $n = 2,125$ ), and “stk” = stacking ( $n = 26,973$ ).

Frame	cP ( $n = 2,334$ )		ncP ( $n = 2,125$ )		stk ( $n = 26,973$ )	
	$R^{\circ}$	$t \text{ \AA}$	$R^{\circ}$	$t \text{ \AA}$	$R^{\circ}$	$t \text{ \AA}$
Sugar-GS	13.35	1.70	19.03	2.48	15.27	1.36
Sugar-Inertial-6	12.32	1.52	18.71	2.42	14.70	1.33
Sugar-Ring	11.85	1.31	18.72	2.30	14.58	1.25
Base-Plane	7.78	0.61	16.42	1.85	12.67	1.07
Base-Center	<b>7.76</b>	<b>0.47</b>	<b>16.41</b>	<b>1.79</b>	<b>12.67</b>	<b>1.06</b>

Canonical cWW pairs strongly favor base-anchored frames in both rotation and translation channels. Noncanonical pairs favor base-anchored frames at the aggregate level, but remain key-heterogeneous: among 191 high-coverage noncanonical relation keys, 76 keys exhibit a sugar/backbone-like rotation pattern with a base-like translation pattern, while by evaluation weight 74% of pairs are base-like in both channels. Stacking is the most heterogeneous: by evaluation weight, 93.6% of stacking pairs favor base-like rotation stability, while 44.6% of stacking translations remain sugar/backbone-like; by key count, 72.7% of stacking keys are base-like in rotation. The multi-state evidence therefore identifies canonical cWW pairing as cleanly base-governed and stacking as a relation in which orientation and translation behave differently.

Table Table 8 reports, for each category, the fraction of keys whose lowest drift comes from a base-family frame. For canonical cWW pairing, a base-family frame achieves the lowest drift for every key in both channels, consistent with the physical organization of Watson-Crick and wobble pairing at the nucleobase level. Noncanonical pairing and stacking are more heterogeneous: only 38% of noncanonical keys have their lowest orientation drift from a base-family frame, yet 68% have their lowest translation drift from a base-family frame; stacking shows 73% of keys favoring a base-family frame in orientation and 68% in translation. Thus, no single frame family yields the lowest drift in both channels for every relation type.

### B.5. Reconstruction Path: Full Single-Frame and Dual-Frame Sweep

The main text reports the core reconstruction tradeoff. Here we report the full reconstruction-path sweep on 51,013 interaction pairs for the gate-passing single frames supported by this path and the three dual configurations that pair a base-anchored frame with the Gram-Schmidt Sugar-GS frame. Sugar-Inertial-4 is excluded from this formal table because it failed the screening gate in Table Table 4.

Among single-frame choices, base-anchored frames preserve the relation geometry identified in the preceding sections but suffer from poor local sugar/backbone reconstruction; sugar/backbone frames give better local reconstruction but cannot accurately recover base atoms. The dual configurations decouple these roles: replacing the global frame with a base-anchored channel substantially reduces base-channel RMSD while retaining the Sugar-GS local reconstruction path.

Overall, the three analyses converge on the same conclusion: selecting any one single frame as a universal representative is inadequate for RNA. Base-anchored frames best capture inter-residue relation geometry (§B.3, §B.4) but do not by themselves solve local sugar/backbone reconstruction (§B.5); sugar/backbone-anchored frames show the complementary behavior. The dual-frame design resolves this by assigning each role to a dedicated channel: a global base frame for relation-bearing geometry and a sugar frame for reconstruction-bearing geometry (§4.1).

Table 8. Fraction of high-coverage relation keys whose lowest drift is achieved by a base-family frame.

Relation category	Keys	Orientation	Translation
Canonical cWW pairs	6	100%	100%
Noncanonical pairs	191	38%	68%
Stacking	139	73%	68%

Table 9. Reconstruction-path comparison on 51,013 interaction pairs.  $\text{RMSD}_{\text{base}}$  measures atoms reconstructed in the base channel;  $\text{RMSD}_{\text{bb}}$  and  $\text{RMSD}_{\text{bridge}}$  measure atoms reconstructed in the sugar/backbone channel and the next-residue  $O3' \rightarrow P$  bridge.

Configuration	$\text{RMSD}_{\text{base}}$ (Å) ↓	$\text{RMSD}_{\text{bb}}$ (Å) ↓	$\text{RMSD}_{\text{bridge}}$ (Å) ↓	Clash rate ↓
<i>Single-frame configurations</i>				
Sugar-GS	6.21	1.88	2.84	$4.9 \times 10^{-4}$
Sugar-Ring	5.37	1.85	2.75	$5.5 \times 10^{-2}$
Sugar-Inertial-6	7.12	1.96	2.30	$8.3 \times 10^{-4}$
Base-Plane	5.85	2.98	4.06	$4.9 \times 10^{-4}$
Base-Inertial	5.60	2.34	2.91	$2.1 \times 10^{-3}$
Base-Center	5.61	3.20	3.74	$2.4 \times 10^{-3}$
<i>Dual-frame configurations with Sugar-GS local frame</i>				
Base-Plane + Sugar-GS	1.58	1.88	2.84	$4.8 \times 10^{-4}$
Base-Inertial + Sugar-GS	1.52	1.88	2.84	$4.8 \times 10^{-4}$
Base-Center + Sugar-GS	1.09	1.88	2.84	$4.9 \times 10^{-4}$

## C. Experiment Details and Ablation Studies

### C.1. Auxiliary Losses, Ablations, and Hyperparameters

The auxiliary losses are defined in Appendix A.2. This appendix reports their empirical ablation and the hyperparameter ranges used in the paper.

#### C.1.1. AUXILIARY OBJECTIVE ABLATION

We ablate auxiliary objectives beyond the core frame, sequence, torsion, and atom23 reconstruction losses. Each variant uses the same architecture, the same 120K training budget, and the same fixed sampler ( $N_T = 50, \alpha = 10$ ). The auxiliary losses are the soft base-template loss, chain-continuity loss, and steric clash loss defined above.

Table 10. Auxiliary objective ablation at fixed sampling ( $N_T = 50, \alpha = 10$ ). The listed weights correspond to auxiliary terms only; all rows keep the core training objective unchanged and use 600 generated structures.

Auxiliary setting	$\lambda_{\text{soft}}$	$\lambda_{\text{chain}}$	$\lambda_{\text{clash}}$	IF scTM ↑	IF scRMSD ↑	GS scTM ↑	GS scRMSD ↑
none	0.0	0.0	0.0	<b>34.33</b>	28.83	<b>26.00</b>	22.33
soft	0.5	0.0	0.0	25.67	23.50	20.17	19.17
soft + chain	0.5	1.0	0.0	32.33	27.83	23.17	22.17
soft + chain + clash	0.5	1.0	5.0	33.50	<b>32.33</b>	25.67	<b>27.17</b>

At the fixed sampler, the core objective without auxiliary terms gives the best scTM-validity under both IF and GS. The soft base-template loss alone substantially reduces validity, suggesting that this auxiliary can over-constrain the sequence-structure channel when used without additional geometric regularization. Adding the chain-continuity term partially recovers performance, and adding the steric clash term gives the best scRMSD-validity, consistent with its role as an atom-level geometric regularizer.

#### C.1.2. HYPERPARAMETER RANGES

Table Table 11 lists the objective weights that appear in the main method, together with the auxiliary-loss ranges used or reserved for ablations.

Table 11. Core objective weights and auxiliary-loss ablation ranges.

Symbol	Reported value(s)	Ablation range	Role
$\tau$	0.25	[0.1, 0.5]	time gate for decoded-geometry supervision
$\lambda_{\text{Seq}}$	1.0	[0.5, 2.0]	nucleotide cross-entropy
$\lambda_{\text{Loc}}$	1.0	[0.25, 2.0]	torsion and decoded atom23 reconstruction objective
$\lambda_{\text{soft}}$	0 or 0.5	[0, 1.0]	soft base-template auxiliary
$\lambda_{\text{chain}}$	0 or 1.0	[0, 2.0]	adjacent $O3'-P$ bridge continuity
$\lambda_{\text{clash}}$	0	[0, 1.0]	optional steric clash auxiliary

Table 12. Augmented Table 2 with Wilson 95% CIs and significance vs. ( $N_T=100, \alpha=20$ ).

$N_T$	$\alpha$	IF		GS	
		Val. scTM % [95% CI]	vs. ref.	Val. scTM % [95% CI]	vs. ref.
50	10	32.00 [28.39, 35.84]	**	24.67 [21.39, 28.27]	*
<b>100</b>	<b>20</b>	<b>40.67 [36.81, 44.64]</b>	—	<b>31.17 [27.59, 34.98]</b>	—
200	10	37.83 [34.04, 41.78]	ns	27.00 [23.60, 30.69]	ns
200	20	38.00 [34.20, 41.95]	ns	28.50 [25.03, 32.24]	ns
300	10	36.00 [32.26, 39.92]	ns	27.50 [24.08, 31.21]	ns
300	20	38.17 [34.37, 42.12]	ns	28.83 [25.35, 32.58]	ns

### C.1.3. SIGNIFICANCE ANALYSIS

All validity metrics are binomial proportions over  $n = 600$  generated structures. We report Wilson 95% confidence intervals (CI) and assess pairwise differences with a two-proportion  $z$ -test. Significance against the reference row of each table is marked as  $*p < 0.05$ ,  $**p < 0.01$ ,  $***p < 0.001$ ; ns:  $p \geq 0.05$ .

Table 12: only (50, 10) differs significantly from the chosen (100, 20); higher budgets show diminishing returns. Table 13: 190K/200K are significantly worse than 180K, supporting the post-180K regression.

880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934

Table 13. Augmented Table 3 with Wilson 95% CIs and significance vs. the 180K checkpoint.

Steps	IF Val. scTM %		GS Val. scTM %	
	Value [95% CI]	vs. 180K	Value [95% CI]	vs. 180K
120K	40.67 [36.81, 44.64]	ns	31.17 [27.59, 34.98]	*
<b>180K</b>	<b>44.00 [40.08, 48.00]</b>	—	<b>37.83 [34.04, 41.78]</b>	—
190K	38.33 [34.53, 42.29]	*	29.17 [25.67, 32.93]	**
200K	36.67 [32.91, 40.60]	**	31.00 [27.43, 34.81]	*